

# upmendexで作る多言語索引 Multilingual index processing by upmendex

田中 琢爾  
TANAKA Takuji

2022年11月19日

# Overview

- Feature of **upmendex**
  - Multilingual index processor
- Localization
  - Latin, Cyrillic, Greek
  - CJK (Chinese, Japanese, Korean)
  - Devanagari, Thai
  - Arabic, Hebrew
- Multilingual environment
- Benchmark

## Index

### — Symbols —

\$ .....	1
¥ .....	1
€ .....	1
2.71828182 .....	1
3.14159265 .....	1

### — C —

Ciudad de México .....	1
------------------------	---

### — I —

İstanbul .....	1
----------------	---

### — S —

São Paulo .....	1
-----------------	---

### — U —

upmendex .....	1
—のインストール .....	1
—の使い方 .....	1
—応用編 .....	1
—入門編 .....	1

### — A —

Aθíva .....	1
-------------	---

### — と —

東京 .....	1
----------	---

### — だ —

대구(大邱) .....	1
--------------	---

대전(大田) .....	1
--------------	---

### — サ —

서울 .....	1
----------	---

### — ハ —

평양(平壤) .....	1
--------------	---

### — ヒ部 —

北京 .....	1
----------	---

### — ハ部 —

廈門(厦门) .....	1
--------------	---

### — ツ部 —

臺北(台北) .....	1
--------------	---

### — カ —

কোলকাতা .....	1
---------------	---

### — ド —

দিল্লী .....	1
--------------	---

# Feature of upmendex (Ver. 1.06)

- Index processor
  - Upper compatible with **MakeIndex/Mendex**
  - Work with upLaTeX/LuaLaTeX/XeLaTeX
- Localization
  - Support 60 Languages / 12 Scripts
    - Latin (incl. non-English), Cyrillic, Greek
    - CJK (Chinese, Japanese, Korean)
    - Devanagari, Thai
    - Arabic, Hebrew
    - Symbol, Number
- Multilingualization
  - Unicode, UTF-8
  - ICU<sup>†</sup> (Collation, Case Conversion, Category Property)
  - Environment for babel/polyglossia

† ICU: International Components for Unicode

# Language, Script, Locale

Language	Script / Index	ICU locale
English	Latin	root
Spanish	Latin	es
German	Latin	de
Turkish	Latin	tr
...		
Russian	Cyrillic	ru
Ukrainian	Cyrillic	uk
...		
Greek	Greek	el
Chinese	Hanzi / Pinyin	zh
	Hanzi / Stroke	zh-u-co-stroke
	Hanzi / Radical-Stroke	zh-u-co-unihan
	Hanzi / Zhuyin	zh-u-co-zhuyin
Japanese	Kana & Hanzi / Kana	ja
Korean	Hangul	ko

Language	Script / Index	ICU locale
Hindi	Devanagari	hi
Marathi	Devanagari	mr
...		
Thai	Thai	th
Persian	Arabic	fa
Arabic	Arabic	ar
...		
Hebrew	Hebrew	he
Yiddish	Hebrew	yi
Common	Symbol Number	

**upmendex** supports 60 languages,  
12 scripts & 95 locales.

# Latin, Cyrillic, Greek

- Sorting (Collation) | ソート順
- Diacritical mark | ダイアクリティカルマーク
- Digraph/Trigraph | ダイグラフ/トライグラフ

# German, Phonebook Sort Order

Collation Rule locale: de-u-co-phonebk

```
&AE<<ä<<<Ä
&OE<<ö<<<Ö
&UE<<ü<<<Ü
...
```

Style File \*.ist locale: de-u-co-phonebk

`icu_locale "de-u-co-phonebk"`

German Inputs in \*.tex

```
ad\index{ad}.
ae\index{ae}.
AE\index{AE}.
ä\index{ä}.
Ä\index{Ä}.
af\index{af}.
...
```

## Index

### — A —

a .....	1
A .....	1
ä .....	1
Ä .....	1
ad .....	1
AD .....	1
ae .....	1
AE .....	1
af .....	1
AF .....	1

default

## Index

### — A —

a .....	1
A .....	1
ad .....	1
AD .....	1
ae .....	1
AE .....	1
ä .....	1
Ä .....	1
af .....	1
AF .....	1

phonebook sort order

# Lithuanian, Sort Order of Y

Collation Rule locale: lt

&I<<j<<<I<<y<<<Y

...

Lithuanian Inputs in \*.tex

i\index{i} .  
I\index{I} .  
j\index{j} .  
J\index{J} .  
y\index{y} .  
Y\index{Y} .  
...

## Rodyklė

### — H —

h ..... 1  
H ..... 1

### — I —

i ..... 1  
I ..... 1  
j ..... 1  
J ..... 1  
y ..... 1  
Y ..... 1

### — J —

j ..... 1  
J ..... 1

### — X —

x ..... 1

### — Z —

z ..... 1

# Slovak, Diacritical Mark

## Collation Rule locale: sk

&O<ô<<<Ô

...

## Collation Rule locale: sk-u-co-search

&L<Í<<<Í<I<<<I  
&O<ó<<<Ó<ô<<<Ô

...

## Slovak Inputs in \*.tex

```
\index{l}.\n\index{í}.\n\index{í}.\n\index{o}.\n\index{ó}.\n\index{ô}.
```

## Index

### — L —

l	.....	1
í	.....	1
I	.....	1

### — O —

o	.....	1
ó	.....	1

### — ô —

ô	.....	1
---	-------	---

default

## Index

### — L —

l	.....	1
---	-------	---

### — Í —

í	.....	1
---	-------	---

### — I —

I	.....	1
---	-------	---

### — O —

o	.....	1
---	-------	---

### — ó —

ó	.....	1
---	-------	---

### — ô —

ô	.....	1
---	-------	---

General-Purpose  
Search

# Turkish, Dotless / Dotted I

language	upper	lower
Turkish	I	i
English	I	i

Collation Rule locale: tr

```
&[before 1]i<ı<<<I
&i<<<İ
...

```

Turkish Inputs in \*.tex

```
h\index{h}.
H\index{H}.
i\index{i}.
ı\index{ı}.
I\index{I}.
İ\index{İ}.
j\index{j}.
J\index{J}.
...

```

Dizin

— H —	
h .....	1
H .....	1
— I —	
i .....	1
ı .....	1
— İ —	
ı .....	1
İ .....	1
— J —	
j .....	1
J .....	1

Turkish

# Hungarian, Digraphs and Trigraph

## Collation Rule locale: hu

```
&C<cs<<<Cs<<<CS
&D<dz<<<Dz<<<DZ
&DZ<dzs<<<Dzs<<<DZS
...
...
```

## Hungarian Inputs in \*.tex

```
cr\index{cr}.
cs\index{cs}.
ct\index{ct}.
dy\index{dy}.
dz\index{dz}.
dzc\index{dzc}.
dzs\index{dzs}.
dzt\index{dzt}.
e\index{e}.
...
```

## Tárgymutató

— C —	
cr .....	1
cr .....	1
ct .....	1
...	
— Cs —	
cs .....	1
...	
— D —	
dy .....	1
...	
— Dz —	
dz .....	1
dzc .....	1
dzt .....	1
dzs .....	1
...	
— Dzs —	
dzs .....	1
...	
— E —	
e .....	1

# Cyrillic & Greek

## Russian Inputs in \*.tex

цветок\index{цветок}.  
птица\index{птица}.  
ветер\index{ветер}.  
луна\index{луна}.

## Greek Inputs in \*.tex

λουλούδι\index{λουλούδι}.  
πουλί\index{πουλί}.  
άνεμος\index{άνεμος}.  
φεγγάρι\index{φεγγάρι}.

## Предметный указатель

— В —	
ветер .....	1
— Л —	
луна .....	1
— П —	
птица .....	1
— Ц —	
цветок.....	1

Russian, Cyrillic

## Ευρετήριο

— Α —	
άνεμος.....	1
— Λ —	
λουλούδι.....	1
— Π —	
πουλί.....	1
— Φ —	
φεγγάρι .....	1

Greek

# CJK (Chinese, Japanese, Korean)

- Chinese: 4 kinds of sort order | 中国語: 4種のソート順
- Japanese: Reading, Extended Kana  
| 日本語: 読み、仮名拡張
- Korean: composed/decomposed | 韓国語: 完成型・組合型

# Chinese, Han Ideograph Sort Order

Style File \*.ist locale: zh-u-co-unihan

```
%icu_locale "zh"
%icu_locale "zh-u-co-stroke"

```

Chinese Inputs in \*.tex

```
花\index{花 (8, 花, huā, ハウル)}
鳥\index{鳥 (11, 鳥, niǎo, ニーム)}
風\index{風 (9, 風, fēng, フン)}
月\index{月 (4, 月, yuè, ユク)}
```

sort order	locale
Pinyin	zh
Stroke	zh-u-co-stroke
Radical-Stroke	zh-u-co-unihan
Zhuyin (Bopomofo)	zh-u-co-zhuyin

# Chinese, Han Ideograph Sort Order

## Pinyin Sort Order 拼音 locale: zh

```
\centerline{\bfseries --- F ---}\par\nobreak
\item 風 (9, 風, fēng, ㄉㄥ) \leaders\hbox{ }\hfill 1
\indexspace

\centerline{\bfseries --- H ---}\par\nobreak
\item 花 (8, 花, huā, ㄏㄨㄚ) \leaders\hbox{ }\hfill 1
```

## Stroke Sort Order 筆畫數 zh-u-co-stroke

```
\centerline{\bfseries --- 四畫 ---}\par\nobreak
\item 月 (4, 月, yuè, ㄩㄝˋ) \leaders\hbox{ }\hfill 1
\indexspace

\centerline{\bfseries --- 八畫 ---}\par\nobreak
\item 花 (8, 花, huā, ㄏㄨㄚ) \leaders\hbox{ }\hfill 1
```

## Radical-Stroke Sort Order 部首筆畫數 zh-u-co-unihan

```
\centerline{\bfseries --- 月部 ---}\par\nobreak
\item 月 (4, 月, yuè, ㄩㄝˋ) \leaders\hbox{ }\hfill 1
\indexspace

\centerline{\bfseries --- 艸部 ---}\par\nobreak
\item 花 (8, 花, huā, ㄏㄨㄚ) \leaders\hbox{ }\hfill 1
```

## Zhuyin (Bopomofo) Sort Order 注音符號 zh-u-co-zhuyin

```
\centerline{\bfseries --- ㄉ ---}\par\nobreak
\item 風 (9, 風, fēng, ㄉㄥ) \leaders\hbox{ }\hfill 1
\indexspace

\centerline{\bfseries --- ㄉ一ㄠˇ ---}\par\nobreak
\item 鳥 (11, 鳥, niǎo, ㄉㄧㄉㄤˇ) \leaders\hbox{ }\hfill 1
```

## upmendex Output \*.ind

# Chinese, Han Ideograph Sort Order

## 索引

<b>- F -</b>	
風 (9, 風, fēng, ㄈㄥ)	.....1
<b>- H -</b>	
花 (8, 花, huā, ㄏㄨㄚ)	.....1
<b>- N -</b>	
鳥 (11, 鳥, niǎo, ㄋㄧㄉㄩㄞˇ)	.....1
<b>- Y -</b>	
月 (4, 月, yuè, ㄩㄝˋ)	.....1

## 索引

<b>- 四畫 -</b>	
月 (4, 月, yuè, ㄩㄝˋ)	.....1
<b>- 八畫 -</b>	
花 (8, 花, huā, ㄏㄨㄚ)	.....1
<b>- 九畫 -</b>	
風 (9, 風, fēng, ㄈㄥ)	.....1
<b>- 十一畫 -</b>	
鳥 (11, 鳥, niǎo, ㄋㄧㄉㄩㄞˇ)	.....1

## 索引

<b>- 月部 -</b>	
月 (4, 月, yuè, ㄩㄝˋ)	.....1
<b>- 艹部 -</b>	
花 (8, 花, huā, ㄏㄨㄚ)	.....1
<b>- 風部 -</b>	
風 (9, 風, fēng, ㄈㄥ)	.....1
<b>- 鳥部 -</b>	
鳥 (11, 鳥, niǎo, ㄋㄧㄉㄩㄞˇ)	.....1

## 索引

<b>- ㄈ -</b>	
風 (9, 風, fēng, ㄈㄥ)	.....1
<b>- ㄋ -</b>	
鳥 (11, 鳥, niǎo, ㄋㄧㄉㄩㄞˇ)	.....1
<b>- ㄏ -</b>	
花 (8, 花, huā, ㄏㄨㄚ)	.....1
<b>- ㄩ -</b>	
月 (4, 月, yuè, ㄩㄝˋ)	.....1

拼音  
pinyin

筆畫數  
stroke

部首筆畫數  
radical-stroke

注音符號  
zhuyin (bopomofo)

# Japanese, Sort by Reading (Yomi)

## Japanese Inputs with Reading (Yomi) in \*.tex

```
\newcommand{\YomiTag}[1]{\relax}
% \index{reading@index_word}
生酒\index{なまざけ@生酒}。
生一本\index{きいっぽん@生一本}。
生け簀\index{いけす@生け簀}。
生絹\index{きぎぬ@生絹}\YomiTag{きぎぬ}%
\index{すずし@生絹}\YomiTag{すずし}。
生飯\index{さば@生飯}。
生姜\index{しょうが@生姜}。
生活\index{せいかつ@生活}\YomiTag{せいかつ}%
\index{たつき@生活}\YomiTag{たつき}。
...
```

Japanese words consist of Hanzi (ideographs) & Kana (syllabaries), sorted by reading (yomi), indexed by Kana.

This feature is implemented by **ASCII mendex**.

## 索引

### — あ —

生憎	1
生け簀	1

### — か —

生一本	1
生絹	1

### — さ —

生飯	1
生姜	1
生絹	1
生活	1

### — た —

生活	1
----	---

### — な —

生酒	1
生業	1

### — は —

生え抜き	1
------	---

# Japanese, Reading & Dictionary

## Japanese Inputs in \*.tex

```
生酒\index{生酒}。
生一本\index{生一本}。
生け簀\index{生け簀}。
生絹\index{生絹}%
  \index{すずし@\生絹\YomiTag{すずし}}。
生え抜き\index{生え抜き}。
...
```

## Dictionary \*.dic

index_word	reading
生酒	なまざけ
生一本	きいっぽん
生け簀	いけす
生絹	きぎぬ
生え抜き	はえぬき
...	

Implemented by **ASCII mendex**.

## 索引

—あ—	—す—
生憎	生絹
.....	.....
1	1
—い—	—せ—
生け簀	生活
.....	.....
1	1
—う—	—た—
生毛	生活
.....	.....
1	1
生まれ付き	
.....	
1	
—な—	
—お—	生きぬ仲
生い立ち	生酒
.....	.....
1	1
—業	生業
—き—	
生一本	—は—
.....	
1	
生絹	生え抜き
.....	.....
1	1
—さ—	—む—
生飯	生す
.....	.....
1	1
—し—	
生姜	
.....	
1	

# Japanese, Hentaigana

## Inputs with Hentaigana in \*.tex

```
後持む\index{後持む}
うふぎ\index{うふぎ}
ゑるお\index{ゑるお}
喜し\index{喜し}
天姫羅\index{天姫羅}
...
```

## Dictionary for Hentaigana \*.dic

index_word	reading
後持む	きそはなしこ
うふぎ	
ゑるお	
喜し	
天姫羅	
...	

## お品書き

— う —		— た —	
うどん	1	だんご	1
うせん	1	巻んぱく	1
うなぎ	1		— て —
うふぎ	1	てんぷら	1
— き —		天姫羅	1
きそば	1		
後持む	1		
— し —			
しるこ	1		
ゑるお	1		
— す —			
すし	1		
喜し	1		
— せ —			
せんべい	1		
せん巻き	1		

# Extended Kana in JIS X 0213

## Extended Kana Inputs in \*.tex

```
が\index{ガ}. % NGA
ぎ\index{ギ}. % NGI
ぐ\index{グ}. % NGU
ぢ\index{ぢ}. % Hiragana Digraph Yori
𠂊\index{𠂊}. % Katakana Digraph Koto
```

## Aynu itak Inputs in \*.tex

```
ク\index{ク}
シ\index{シ}
ス\index{ス}
ブ\index{ブ}
ゼ\index{ゼ}
ヅ\index{ヅ}
ド\index{ド}
...
```

## 索引

### 一か一

か	1
が	1
が	1
ガ	1
ギ	1
ギ	1
グ	1
げ	1
ゲ	1
ご	1
ゴ	1
𠂊	1

### 一や一

ぢ	1
---	---

## 索引

### 一か一

ク	1
ク	1

### 一さ一

シ	1
ス	1
ゼ	1

### 一た一

ヅ	1
ト	1
ド	1

### 一な一

ヌ	1
---	---

### 一は一

ブ	1
ブ	1

# Japanese, Archaic Kana



Katakana “YE” = 「ヱ」 ???  
“WE” ≠ 「ヱ」 ???

# Japanese, Archaic Kana



Katakana “YE” = 「ヱ」 ???  
“WE” ≠ 「ヱ」 ???

A page from a Japanese kana reader titled "1997.28.1445". It contains several horizontal rows of Japanese text. The first row includes the character 'ヱ'. The text is arranged in columns, likely illustrating the use of the character 'ヱ' in different contexts or readings.

ンワラヤマハ	ナタサカア	五
ヰリレミヒニチシキイ		三
于ルユムフヌツスクウ		音
エレエメヘネテセケエ		四
ヲロヨモホノトソコホ		五

Katakana Letter Archaic YE = 「ヱ」  
defined in Unicode 14.0 (2021)

Ref. <http://www.unicode.org/L2/L2019/19381-missing-kana.pdf>  
<https://dl.ndl.go.jp/info:ndljp/pid/993592>

# Japanese, Archaic Kana

## Archaic Kana Inputs in \*.tex

```
\text{I}\index{\text{I}}. % Hiragana Archaic YE
% or Hentaigana E-1
\text{I}\index{\text{I}}. % Hiragana WU
\text{I}\index{\text{I}}. % Katakana YI
\text{I}\index{\text{I}}. % Katakana YE
\text{I}\index{\text{I}}. % Katakana WU
...
```

Dictionary: “I” is Hentaigana E-1

I え

Style: “I” is Hiragana Archaic YE

icu\_rules "&I<\text{I}<<<\text{I}<\text{Y}"

## 索引

— イ —	
え .....	1
工 .....	1
I (Hentaigana E-1) .....	1
— ャ —	
ヤ .....	1
— レ —	
レ .....	1
— ュ —	
ユ .....	1
— イ —	
ヰ .....	1
— ョ —	
ヨ .....	1
— ナ —	
け .....	1
于 .....	1

## 索引

— イ —	
え .....	1
工 .....	1
— ャ —	
ヤ .....	1
— レ —	
レ .....	1
— ュ —	
ユ .....	1
— イ —	
ヰ .....	1
I (Archaic YE) .....	1
— ョ —	
ヨ .....	1
— ナ —	
け .....	1
于 .....	1

# Korean, Modern / Archaic Hangul



	Unicode block	style	upmendex	upLaTeX	XeLaTeX
modarn ex. 일	Hangul Syllables	composed	완성형	✓	✓
	Hangul Jamo	decomposed	조합형	✓	N.A.
archaic ex. ·실	Private Use Area	composed	완성형	via dictionary	✓
	Hangul Jamo	decomposed	조합형	✓	N.A.

# Korean Hangul

## Hangul Inputs in \*.tex

```

쓰\index{쓰 (composed)}.
々-\index{々- (decomposed)}.

々 .\index{々 . (archaic)}.
々 . :\index{々 . : (archaic with tone mark)}.
々\index{々 (Hanyang PUA)}.

...

```

## Dictionary for PUA code \*.dic

### Hanyang PUA code decomposed

쓰	々 .
쓰	々 .
...	

## 찾아보기

— ^ —	
스 (composed) .....	1
々 (archaic) .....	1
— ^ —	
쓰 (composed) .....	1
쓰 (decomposed) .....	1
:々 (archaic with tone mark) .....	1
々 (archaic) .....	1
々 (Hanyang PUA) .....	1
— ^ —	
々 (archaic) .....	1
々 (Hanyang PUA) .....	1
— ^ —	
쓰 (archaic) .....	1
— ^ —	
쓰 (archaic) .....	1
— ^ —	
쓰 (archaic) .....	1

# Complex Text Layout

- Devanagari, Thai
- Arabic, Hebrew: R-to-L typeset
- Symbol, Number

# Devanagari & Thai (experimental)

## Hindi Inputs in \*.tex

```
फूल \index{फूल}
चिड़िया \index{चिड़िया}
हवा \index{हवा}
चांद \index{चांद}
...

```

## Thai Inputs in \*.tex

```
ດອກໄມ້ \index{ດອກໄມ້}
ນກ \index{ນກ}
ລມ \index{ລມ}
ດວງຈັນທີ່ \index{ດວງຈັນທີ່}
...

```

### सूची

---	च ---
चांद	..... 1
चिड़िया	..... 1
---	फ ---
फूल	..... 1
---	ह ---
हवा	..... 1

### ດຽວចະຫິນ

---	ອ ---
ດວງຈັນທີ່	..... 1
ດອກໄມ້	..... 1
---	ຸ ---
ນກ	..... 1
---	າ ---
ລມ	..... 1

Hindi, Devanagari

Thai

Typeset by XeLaTeX

# Arabic & Hebrew (experimental)

## Arabic Inputs in \*.tex

```
\index{زهرة}
\index{عصفوري}
\index{ريح}
\index{القمر}
```

...

## Hebrew Inputs in \*.tex

```
\index{פָנָח}
\index{צִיפּוֹר}
\index{רוֹעַ}
\index{ירָח}
```

...

الفهرس	
--   --	القمر
-- ر --	ريح
-- ز --	زهرة
-- ع --	عصفوري

Arabic

مفتاح	
-- ' --	يرّاح
-- פ --	פרה
-- צ --	ציפור
-- ר --	רוֹעַ

Hebrew

R-to-L typeset by XeLaTeX.  
**upmendex** processes only indexing.

# Symbol, Number

Script	charType	example	treatment by upmendex
Latin	Lu, Ll, Lo, ... : letters	ABCabc A a Ⓐ	directly pass to ICU collator
Greek	Lu, LL, Lo, ... : letters	ΑΒΓαβγ	direct
Cyrillic	Lu, LL, Lo, ... : letters	АБВабв	direct
Kana	Lo : other letter	あいうアアヲヲト	direct
Hangul	Lo : other letter	가나다Forgeable	direct
Hanzi	Lo : other letter	花鳥風月	lookup dictionary or direct
—	Lm : modifier letter	— ° -	direct
Number	Nd : decimal digit number No : other number	0120 1 2 12③④⑤⑥7,(8)9.	direct lookup dictionary or direct
Symbol	Sk : modifier symbol Sm : math symbol So : other symbol Sc : currency symbol Po, Pd, Mn, Me, ... : other punctuation etc.	÷▷# ◐◑●♥●●●● €\$ \$ ₩¥ ¥ ₩₩ ?!?¡¿† # § ¶ —	direct lookup dictionary or direct lookup dictionary or direct lookup dictionary or direct direct
—	Cc : control character	ESC, BS, DEL	ignore
—	Cf : format character	BOM, RLM	dict
others	Lo, etc. (unknown scripts)		lookup dic or direct (option “-f”) or ignore

Characters are classified by Unicode *General Category Values* or “charTypes”

Ref. [https://unicode.org/reports/tr44/#General\\_Category\\_Values](https://unicode.org/reports/tr44/#General_Category_Values)

# Multilingual Environment with upLaTeX/pxbabel

## Block Setting for Scripts in Style File \*.ist

```

script_preamble cyrillic "\n\\fontencoding{T2A}\\selectfont"
script_postamble cyrillic "\n\\fontencoding{T1}\\selectfont"

script_preamble hangul "\n\\begin{otherlanguage}{korean}"
script_postamble hangul "\n\\end{otherlanguage}"

script_preamble hanzi "\n\\begin{otherlanguage}{tchinese}"
script_postamble hanzi "\n\\end{otherlanguage}"

```

## upmendex Output \*.ind

```

\centerline{\bfseries --- C ---}\par\nobreak
\item София\leaders\hbox{~}\hfill 1
...
\fontencoding{T1}\selectfont

\indexspace

\centerline{--- サ ---}\par\nobreak
\item さいたま\leaders\hbox{~}\hfill {1}
\item 札幌\leaders\hbox{~}\hfill {1}
...
\begin{otherlanguage}{korean}

\indexspace

\centerline{\bfseries --- ハ ---}\par\nobreak
\item 부산(釜山)\leaders\hbox{~}\hfill 1
...
\end{otherlanguage}

```

# Multilingual Environment with XeLaTeX/polyglossia

## Block Setting for Scripts in Style File \*.ist

```
script_preamble cyrillic "\n\\begin{russian}"  
script_postamble cyrillic "\n\\end{russian}"  
  
script_preamble kana "\n\\begin{japanese}"  
script_postamble kana "\n\\end{japanese}"  
  
script_preamble hangul "\n\n\\begin{korean}"  
script_postamble hangul "\n\\end{korean}"  
  
script_preamble hebrew "\n\\begin{hebrew}"  
script_postamble hebrew "\n\\end{hebrew}"
```

## upmendex Output \*.ind

```
\centerline{--- さ ---}\par\nobreak  
  \item さいたま\leaders\hbox{~}\hfill {1}  
...  
\end{japanese}  
  
\begin{korean}  
  \indexspace  
  
\centerline{--- ㄷ ---}\par\nobreak  
  \item 대구(大邱)\leaders\hbox{~}\hfill {1}  
...  
\end{korean}  
  
\begin{hebrew}  
  \indexspace  
  
\centerline{--- א ---}\par\nobreak  
  \item תַּדְבִּיר\leaders\hbox{~}\hfill {2}  
...  
\end{hebrew}
```

# Output of Multilingual Index

## 索引

— symbols —		— サ —
€	.....	1 さいたま.....1
3.14159265	.....	1 札幌.....1
— c —		— ト —
Ciudad de México	.....	1 東京.....1
— i —		— ン —
İstanbul	.....	1 대구(大邱).....1 대전(大田).....1
— s —		— ァ —
São Paulo	.....	1 서울.....1
— ө —		— Ⅱ —
Београд	.....	2 平양(平壤).....1
Бишкек	.....	2 平양(平壤).....1
— қ —		— 五畫 —
Київ	.....	2 北京.....1
— м —		— 十三畫 —
Москва	.....	2 廈門(厦门).....1
— オ —		— 十四畫 —
大阪	.....	1 臺北(台北).....1

with upLaTeX & pxbabel

## Index

— Symbols —		— 파 —
€	.....	1 평양(平壤).....1
3.14159265	.....	1 — 匕部 —
— I —		北京.....1
İstanbul	.....	— 至部 —
— S —		臺北(台北).....1
São Paulo	.....	— 𩚴 —
— A —		दिल्ली.....2
Aθήνα	.....	— 𩚵 —
— Θ —		ముర్బీ.....2
Θεσσαλονίκη	.....	— 𩚷 —
— K —		గ్రూహమహానాద.....2
Київ	.....	— 𩚸 —
— M —		սանկուչ.....2
Москва	.....	— 𩚹 —
— さ —		أبو ظبی.....2
さいたま	.....	— 𩚻 —
— と —		دیوب.....2
東京	.....	— 𩚼 —
— サ —		ירושלים.....2
서울	.....	— 𩚽 —
		תל אביב.....2

with XeLaTeX & polyglossia

# Benchmark

	makeindex	mendex	upmendex	xindy
internal encoding	8bit 1byte	EUC-JP	UTF-16	Unicode
Collator	locale	ASCII, Kana	ICU collator	
Latin	✓	✓ ASCII	Lang:37, Locale:62	Lang:32
Greek			Lang:1, Locale:1	Lang:1
Cyrillic			Lang:9, Locale:9	Lang:6
Chinese			Lang:1, Locale:4 ✓ (Yomi, Dict)	
Japanese		✓ (Yomi, Dict)	Lang:1, Locale:2	
Korean			Lang:1, Locale:3	
Devanagari			Lang:3, Locale:3	
Thai			Lang:1, Locale:1	
Arabic			Lang:6, Locale:7	
Hebrew			Lang:2, Locale:3	Lang:1
Other				Lang:4
Total		Lang:2	Lang:60, Locale:95	Lang:44

# Languages by Number of Native Speakers

	Language	Script	speakers	ICU	Polyglossia	upmendex	xindy
1	Chinese	Hanzi	1,370,000,000	✓		✓	
2	English	Latin	530,000,000	✓	✓	✓	✓
3	Hindi	Devanagari	490,000,000	✓	✓	✓	
4	Spanish	Latin	420,000,000	✓	✓	✓	✓
5	Arabic	Arabic	230,000,000	✓	✓	✓	
6	Bengali	Bengali	220,000,000	✓	✓		
7	Portuguese	Latin	215,000,000	✓	✓	✓	✓
8	Russian	Cyrillic	180,000,000	✓	✓	✓	✓
9	Japanese	Kana & Hanzi	134,000,000	✓	✓	✓	
10	German	Latin	130,000,000	✓	✓	✓	✓
11	French	Latin	123,000,000	✓	✓	✓	✓
12	Punjabi	Gurmukhi	90,000,000	✓			
13	Javanese	Latin	75,000,000	✓	✓	✓	
14	Korean	Hangul	75,000,000	✓	✓	✓	
15	Vietnamese	Latin	70,000,000	✓	✓	✓	✓
16	Telugu	Telugu	70,000,000	✓	✓		
17	Marathi	Devanagari	68,000,000	✓	✓	✓	
18	Tamil	Tamil	74,000,000	✓	✓		
19	Persian	Arabic	46,000,000	✓	✓	✓	
20	Urdu	Arabic	61,000,000	✓	✓	✓	

Ref. <https://ja.wikipedia.org/wiki/ネイティブスピーカーの数が多い言語の一覧>

# Languages used on the Internet

	Language	share	Script	ICU	pg	upm	xnd
1	English	63.4 %	Latin	root	✓	✓	✓
2	Russian	7.1 %	Cyrillic	ru	✓	✓	✓
3	Spanish	3.9 %	Latin	es	✓	✓	✓
4	German	3.7 %	Latin	de	✓	✓	✓
5	Turkish	3.5 %	Latin	tr	✓	✓	✓
6	Persian	2.5 %	Arabic	fa	✓	✓	
7	French	2.0 %	Latin	root	✓	✓	✓
8	Japanese	1.9 %	Kana	ja	✓	✓	
9	Portuguese	1.8 %	Latin	pt	✓	✓	✓
10	Chinese	1.3 %	Hanzi	zh		✓	
11	Vietnamese	1.3 %	Latin	vi	✓	✓	✓
12	Italian	1.0 %	Latin	root	✓	✓	✓
13	Arabic	0.9 %	Arabic	ar	✓	✓	
14	Polish	0.9 %	Latin	pl	✓	✓	✓
15	Greek	0.7 %	Greek	el	✓	✓	✓
16	Dutch	0.7 %	Latin	nl	✓	✓	✓
17	Indonesian	0.7 %	Latin	root	✓	✓	
18	Korean	0.6 %	Hangul	ko	✓	✓	
19	Czech	0.4 %	Latin	cs	✓	✓	✓
20	Thai	0.4 %	Thai	th	✓	✓	

	Language	share	Script	ICU	pg	upm	xnd
21	Ukrainian	0.3 %	Cyrillic	uk	✓	✓	✓
22	Hebrew	0.3 %	Hebrew	he	✓	✓	✓
23	Swedish	0.3 %	Latin	sv	✓	✓	✓
24	Romanian	0.3 %	Latin	ro	✓	✓	✓
25	Hungarian	0.3 %	Latin	hu	✓	✓	✓
26	Danish	0.2 %	Latin	da	✓	✓	✓
27	Slovak	0.2 %	Latin	sk	✓	✓	✓
28	Serbian	0.2 %	Latn, Cyril	sr	✓	✓	✓
29	Bulgarian	0.1 %	Cyrillic	bg	✓	✓	✓
30	Finnish	0.1 %	Latin	fi	✓	✓	✓
31	Croatian	0.1 %	Latin	hr	✓	✓	✓
32	Lithuanian	0.1 %	Latin	lt	✓	✓	✓
33	Norwegian (Bokmål)	0.1 %	Latin	nb	✓	✓	✓
34	Hindi	0.1 %	Devanagari	hi	✓	✓	
35	Norwegian (nynorsk)	0.1 %	Latin	nn	✓	✓	✓
36	Slovenian	0.1 %	Latin	sl	✓	✓	✓
37	Latvian	0.1 %	Latin	lv	✓	✓	✓
38	Estonian	0.1 %	Latin	et	✓	✓	✓
39	Azerbaijani	< 0.1 %	Latin	az		✓	
40	Catalan	< 0.1 %	Latin	root	✓	✓	

Ref. <https://ja.wikipedia.org/wiki/インターネットにおける言語の使用>

# To Do

- Support more scripts
  - Bengali, Telugu, Tamil, Malayalam, Kannada, Gujarati, Oriya, Sinhala
  - Khmer, Lao, Myanmar (Burmese)
  - Tibetan, Mongolian
  - Armenian, Georgian
  - Ethiopic (Amharic)
  - etc.
- Support more locales
  - Latin Script: Sorbian, Hausa, Igbo, Yoruba, Kalaallisut, Breton, Uzbek
  - etc.
- Add style options
  - script\_head

## Feedback is welcome

<https://github.com/t-tk/upmendex-package/issues>

# Summary

I introduced multilingual index processor  
**upmendex**.

- Feature
- Localization: 60 Languages, 12 Scripts
  - Latin, Cyrillic, Greek
  - CJK (Chinese, Japanese, Korean)
  - Devanagari, Thai, Arabic, Hebrew
  - Symbol, Number
- Multilingualization
  - Environment for  
 upLaTeX/babel & XeLaTeX/polyglossia

## Index

— Symbols —	— と —
\$ ..... 1	東京 ..... 1
¥ ..... 1	— 다 —
€ ..... 1	대구(大邱) ..... 1
2.71828182 ..... 1	대전(大田) ..... 1
3.14159265 ..... 1	
— さ —	
— サ —	
Ciudad de México ..... 1	서울 ..... 1
— い —	
İstanbul ..... 1	평양(平壤) ..... 1
— イ —	
São Paulo ..... 1	— ヒ部 —
— ザ —	
upmendex ..... 1	北京 ..... 1
— う —	
— のインストール ..... 1	— フ部 —
— の使い方 ..... 1	廈門(厦门) ..... 1
— 応用編 ..... 1	— 至部 —
— 入門編 ..... 1	臺北(台北) ..... 1
— あ —	
Aθήνα ..... 1	— ク —
— あ —	
	কোলকাতা ..... 1
— あ —	
	— だ —
	দিল্লী ..... 1

# References

- ① ASCII Nihongo TeX (Publishing TeX), ASCII MEDIA WORKS (web site by DWANGO Co., Ltd.).  
The site distributes mendex source files.
- ② Source/Document distribution of upmendex — multilingual index processor @ GitHub.  
upmendex @ CTAN
- ③ upTeX, upLaTeX — unicode version of pTeX, pLaTeX
- ④ International Components for Unicode (ICU)
- ⑤ PXbase — LaTeX: Support library for other PX packages @ GitHub. The repository  
distributes pxbabel.  
pxbase @ CTAN
- ⑥ polyglossia — An alternative to Babel for XeLaTeX and LuaLaTeX @ GitHub.  
polyglossia @ CTAN
- ⑦ “Indexing Makes Your Book Perfect” by SHIKANO Keiichiro at TUG2013, October, Tokyo.